

# Word Representation

## I-hot representation:

$$V = [a, aaron, man, woman, King, Queen, Apple, Orange \dots]$$

Man (5391)   Woman (9853)   King (4914)   Queen (7157)   Apple (456)   Orange (6257)

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$
$O_{5391}$	$O_{9853}$				

Problem: Correlations between any  $\geq 2$  words are 0

In other word, Networks using this type of word embeddings don't understand meanings and relations of words.

e.g. I want a glass of orange juice

I want a glass of apple ?

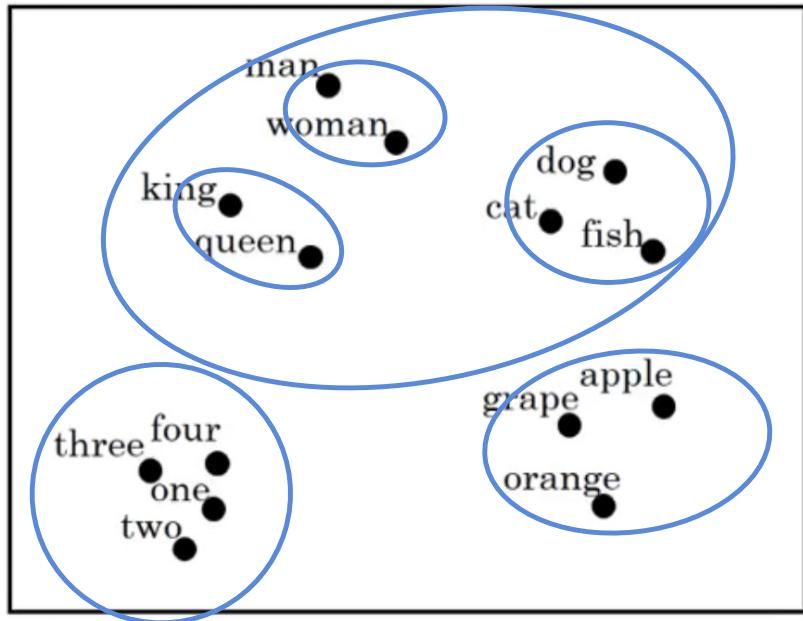
Lack of understanding of similarities between words.

Featurized representation: word embedding

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
size cost altit. verb	$\vdots$	$\vdots$				
	$e_{5391}$	$e_{9853}$				

300D vector

# Visualizing word embeddings



t-SNE =

map 300D feature vector  
to 2D space

- Preserving spatial relationships
- Not a linear mapping  
→ lengths are distorted

## Transfer Learning and word embeddings

Common Practice:

- 1) Learn word embeddings from large corpus (1-100B words)  
(or download pre-trained embedding online)
- 2) Transfer embedding to new task with smaller  
training set (say, 100k words)
- 3) (Optional) Continue to finetune the word embeddings  
with new data.

Word Embedding in NLP is similar to face / image  
encoding in facial recognition task.

# Properties of Word Embeddings

motivation problem :

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

Given the above word embeddings, how can computers solve the problem below?

Man  $\rightarrow$  Woman as King  $\rightarrow$  ?

One possible approach:

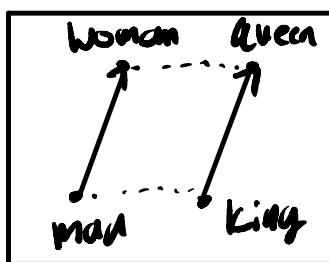
$$e_{\text{man}} - e_{\text{woman}} \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{?}$$

$$e_{\text{king}} - e_{\text{queen}} \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$e_{?}$   
queen

Analogy Using Word Vectors:



300D

$$e_{\text{man}} - e_{\text{woman}} = e_{\text{king}} - e_{?} \approx e_{\text{queen}}$$

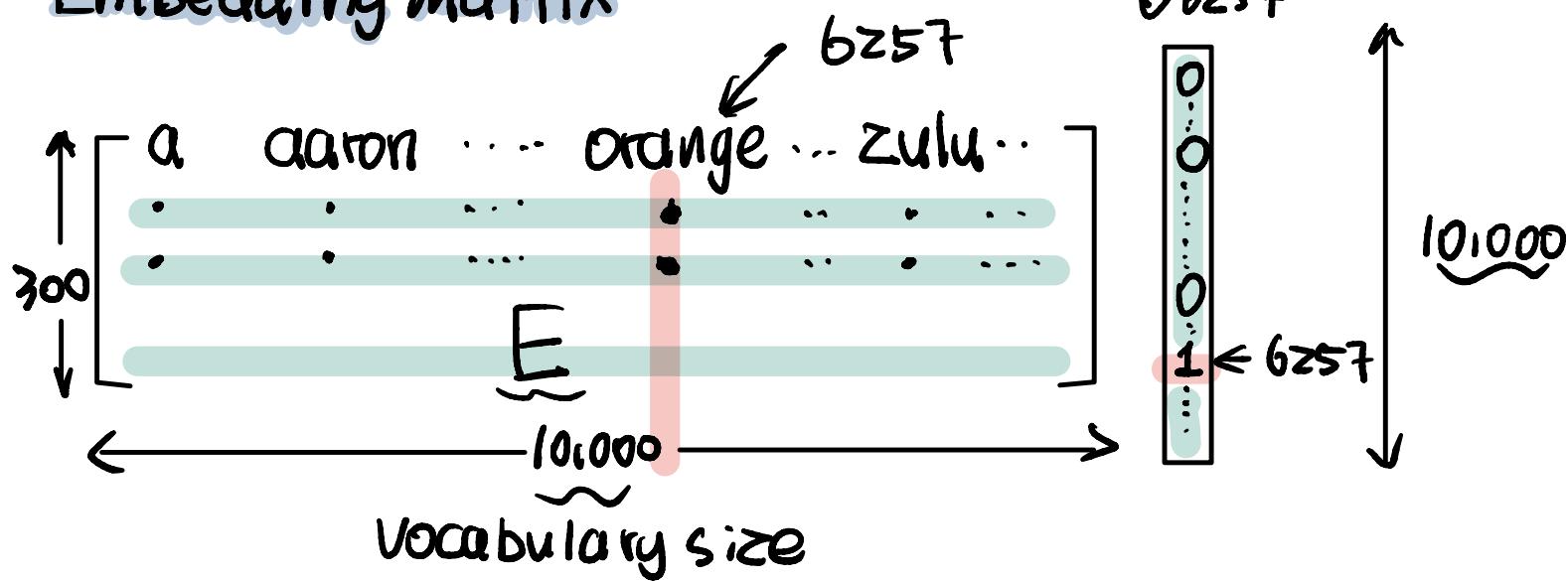
Find word w:

$$\operatorname{argmax}_w \operatorname{Sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

$$\operatorname{Sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

# Learning Word Embeddings

## Embedding Matrix



$$E \cdot \underbrace{O_{6257}}_{\text{Shape} = (300, 1)}$$

→ embedding for  $O_{6257}$  (Orange)

$$\underline{E \cdot O_j = e_j}$$

= embedding for word  $j$

## Learn word embedding

I	want	a	glass	of	orange
4343	9665	1	3852	6163	$e_{4343} = E o_{4343}$

Prediction

I	$o_{4343}$	→	$E$	→	$e_{4343}$
want	$o_{9665}$	→	$E$	→	$e_{9665}$
a	$o_1$	→	$E$	→	$e_1$
glass	$o_{3852}$	→	$E$	→	$e_{3852}$
of	$o_{6163}$	→	$E$	→	$e_{6163}$
orange	$o_{6257}$	→	$E$	→	$e_{6257}$

Neural language

model

Learn embeddings thru training network

softmax  
10,000

$w_2, b_2$

$w_1, b_1$

\*Same embedding matrix for each word

## Other Context/target pairs

I want a glass of orange juice to go along with cereal

A red curved arrow points from the word 'orange' to the word 'context'. A blue vertical arrow points from the word 'juice' to the word 'target'.

Possible context:

- Last 4 word
- Last 1 word
- 4 words on left and right
- Nearby 1 word → skip gram.

...

## word2vec

I want a glass of orange juice to go along with cereal

skip gram:

Context:

Orange  
Orange  
orange

Target:

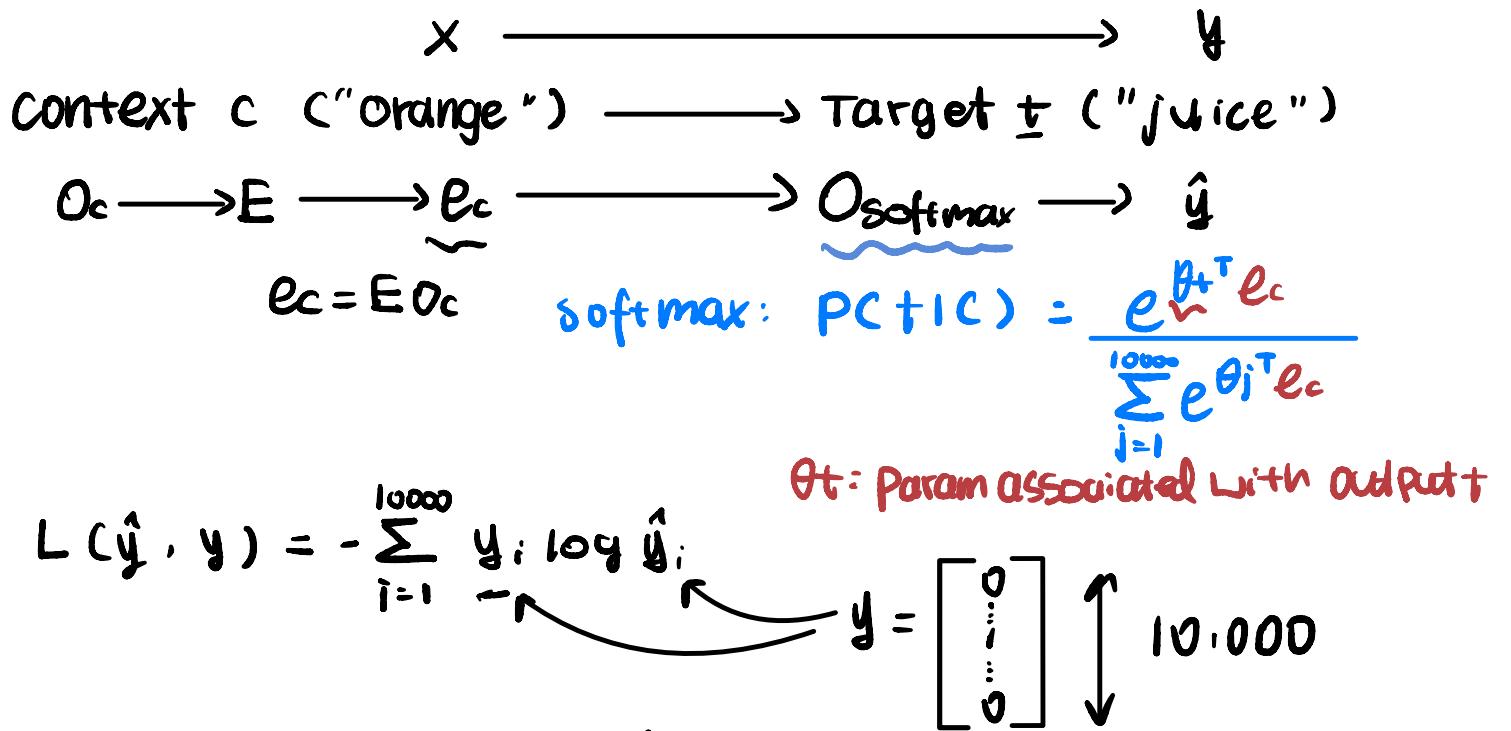
Juice  
glass  
my

a random chosen  
word with a window  
close to context

As we can guess, the training data generated in this way won't do well on the supervised learning task which maps context to target. However, it's proven to be very effective for the network to learn the embedding matrix.

model:

Vocab size = 10,000K



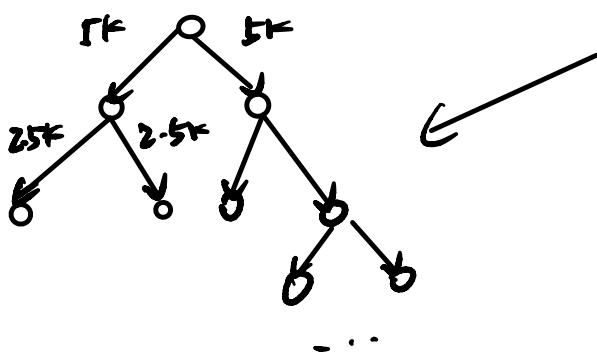
problems with softmax classification

softmax:  $P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$

This computation of this summation over all vocabulary is very expensive

One direct solution:

Hierarchical softmax

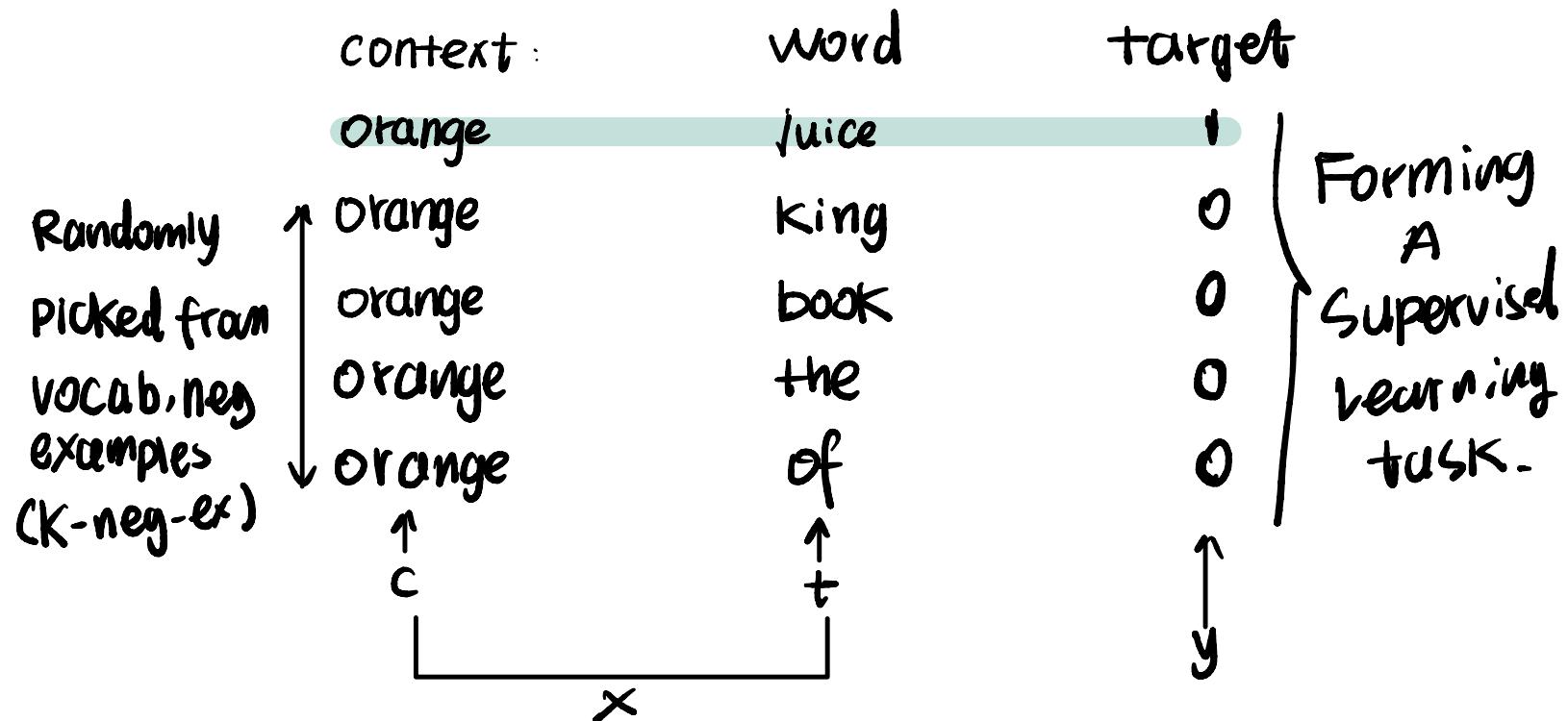


Won't be a balanced binary tree: less frequent words will be buried deeper down the tree.

\* $P(c)$ : choosing context word won't happen within a uniform distribution. Avoid drawing too many common words such as the, of, a, ...

Negative Sampling : Another algorithm that avoids the softmax problem.

I want a glass of orange juice to go along with cereal



$$P(y=1 | c, t) = \sigma(\theta_t^T e_c)$$

Orange  
6257

$$0_{6257} \longrightarrow E \longrightarrow e_{6257}$$

Instead of calculating softmax over entire vocabulary set, we only perform k+1 binary classification problems. (Kneg + 1 pos example).

Selecting words to generate negative examples :  $P(w_i) = \frac{f(w_i)^{3/4}}{\sum_j f(w_j)^{3/4}}$

for juice ?  
for book ?  
10,000 binary classification problem.

for King ?  
for ...

## Glove word vectors

I want a glass of orange juice to go along with cereal

c, t

$X_{ij} = \# \text{ times } j \text{ appears in context of } i$

$\begin{matrix} \uparrow \\ c \end{matrix}$        $\begin{matrix} \uparrow \\ t \end{matrix}$

$\downarrow$

$\begin{matrix} \uparrow \\ c \end{matrix}$

Sometimes  $X_{ij} = X_{ji}$ , depending on how "x appears in context of y" is defined.

Model:

$$\text{minimize } \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(X_{ij}) (\theta_i^T e_j + b_i + b_j - \log X_{ij})^2$$

" $\theta_i^T e_c$ "

Correlation  $\approx \# \text{ times } X_{ij}$

Weighting term:

$$f(X_{ij}) = 0 \text{ if } X_{ij} = 0 \rightarrow 0 \log 0 = 0 \quad \theta_i, e_i \text{ are symmetric}$$

destressing  $\rightarrow$  this, is, a, an ...

$$e_w^{(\text{final})} = \frac{e_w + \theta_w}{\Sigma}$$

addressing  $\rightarrow$  durian ...

\* Featurization might not be interpretable.

$$\rightarrow (A\theta_i)^T (A^{-T} e_j) = \theta_i^T A^T A^{-T} e_j = \theta_i^T e_j$$

Linear transformations may shift featurization away from our natural interpretation.

# Sentiment Classification:

mapping from  $x$   $\longrightarrow$   $y$

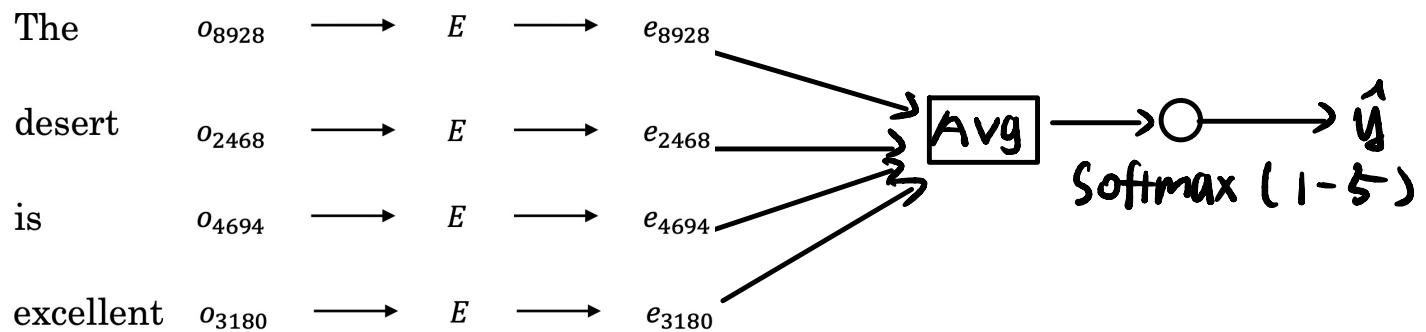
The dessert is excellent.  $\star\star\star\star\star$

Service was quite slow.  $\star\star\star\star\star$

Good for a quick meal, but nothing special.  $\star\star\star\star\star$

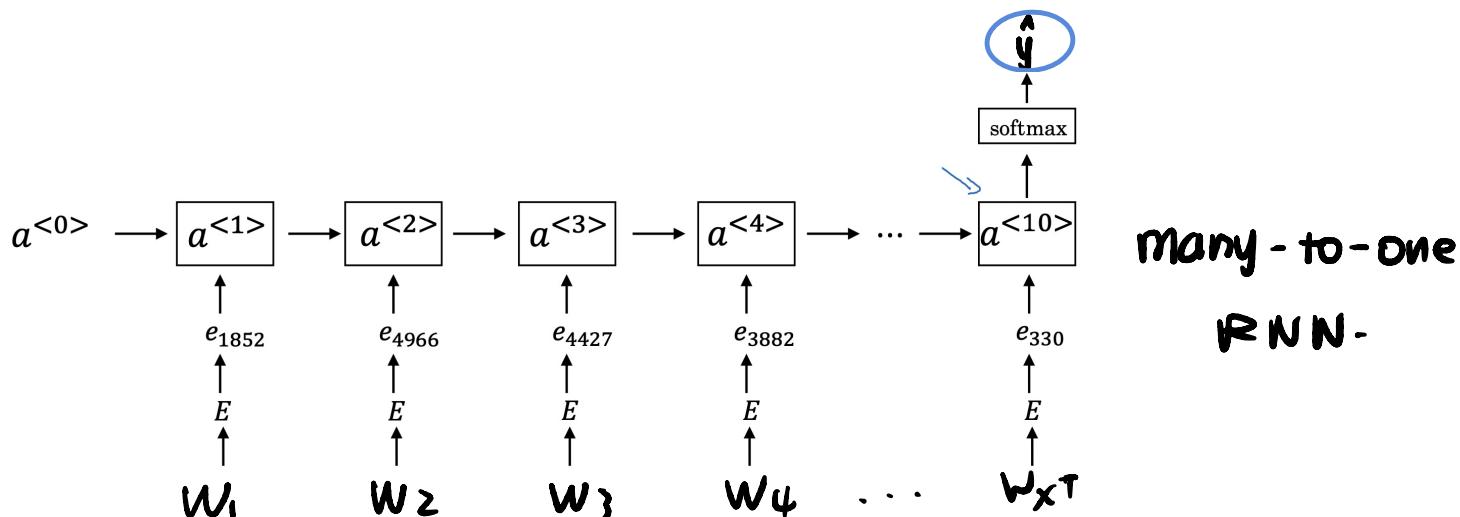
Completely lacking in good taste, good service, and good ambience.  $\star\star\star\star\star$

## Simple sentiment classification model:

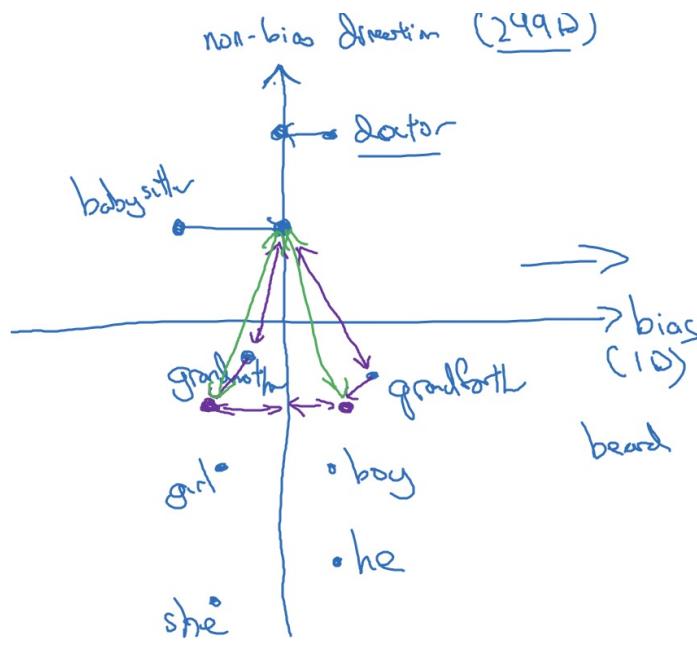


This simple model performs very poorly when input sentence contains more complex logic and sentiments e.g. "Completely lacking in good taste and good service"

## Solution: RNN for Sentiment classification



# Debiasing word embeddings



1. Identify bias direction.

{  
he - she  
male - female  
:  
average}

2. Neutralize: For every word that is not definitional, project to get rid of bias.

3. Equalize pairs.

{  
→ girl - boy }

[Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings] ↵

Andrew Ng